



TITLE:

ベクトル値マルコフ決定過程における最適定常政策(数理計画モデルにおける最適化理論)

AUTHOR(S):

涌田, 和芳

CITATION:

涌田, 和芳. ベクトル値マルコフ決定過程における最適定常政策(数理計画モデルにおける最適化理論). 数理解析研究所講究録 1992, 798: 171-179

ISSUE DATE:

1992-08

URL:

<http://hdl.handle.net/2433/82794>

RIGHT:

ベクトル値マルコフ決定過程における最適定常政策

長岡高専 涌田和芳 (Kazuyoshi Wakuta)

§ 1. 序

割引因子をもつベクトル値マルコフ決定過程において、最適定常政策の特徴付けを考える。Furukawa [1980] は、この問題について最初に議論し、次のような結果を発表している：ある連続性条件の仮定の下で、定常政策が定常政策全体の中で最適であるための必要十分条件は、その利得が最適方程式

$$u(i) \in e \left\{ \bigcup_{a \in A} \sum_{j \in S} (r_{ij}^a + \beta u(j)) g_{ij}^a \right\}$$

の maximal な解となることである。一方、White [1982] は、どんな定常政策によっても支配されない非定常政策の存在を示している。ここでは、すべての (randomized, history-dependent な) 政策の中で最適定常政策の特徴付けることを考える。

§ 2. ベクトル値マルコフ決定過程

ベクトル値マルコフ決定過程 (VMDP) は, 次のもので定義される.

S : 空でないボレル集合, 状態空間

A : " , 行動空間

F : 行動制約集合関数 $F(s) \subset A$ ($s \in S$), $G_r F = \{(s, a) \mid s \in S, a \in F(s)\}$ は, $S \times A$ のボレル集合で, S から A への 1 つのボレル可測写像のグラフを含むとする

$g \in Q(S \mid G_r F)$: システムの運動法則, ただし, $Q(Y \mid X)$ は X から Y への推移確率の全体を表わす

$r \in M^p(G_r F)$: 利得関数, ただし, $M^p(X)$ は X から R^p ($p \geq 1$) への有界ボレル可測関数の全体を表わす

β ($0 \leq \beta < 1$): 割引因子

政策 π の利得は,

$$I(\pi)(s_1) = E_\pi \left[\sum_{i=1}^{\infty} \beta^{i-1} r(s_i, a_i) \mid s_1 \right], \quad s_1 \in S$$

で定義し,

$$V(s_1) = \bigcup_{\pi \in \Pi} \{ I(\pi)(s_1) \}, \quad s_1 \in S$$

とおく. ただし, Π は政策全体の集合を表わす. そして,

$e(V(\Delta_i))$, $\Delta_i \in S$ を, nontrivial, closed convex cone K を domination cone とする $V(\Delta_i)$, $\Delta_i \in S$, の properly efficient な元の集合とする.

定義 $I(\pi^*)(\Delta_i) \in e(V(\Delta_i))$, $\Delta_i \in S$, が成り立つとき, π^* は最適であるという.

§ 3. 最適定常政策

各 $\Delta_i \in S$ に対して $c(\Delta_i) \in R^P$ を選び, 利得関数 $r(\Delta_i, \Delta_n, a_n) = \langle c(\Delta_i), r(\Delta_n, a_n) \rangle$ をもつ非定常動的計画 (NDP(c)) を考える. ただし, $\langle \cdot, \cdot \rangle$ は内積を表わす. その他は, VMDP と同じとする. NDP(c) における政策 π の利得は,

$$J(\pi)(\Delta_i) = E_{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} r(\Delta_i, \Delta_n, a_n) \mid \Delta_i \right], \Delta_i \in S$$

で定義し, $J(\pi^*)(\Delta_i) \geq J(\pi)(\Delta_i)$, $\Delta_i \in S$, $\pi \in \Pi$ が成り立つとき π^* は最適であるという.

仮定 K の dual cone $K^* = \{ k^* \in R^P \mid \langle k^*, k \rangle \geq 0, k \in K \}$ に対して, $\text{int } K^* \neq \emptyset$ と仮定する.

命題 3.1 π^* がある $c(\Delta_i) \in (\text{int } K^*)$, $\Delta_i \in S$, をもつ NDP(c) に対して最適ならば, VMDP に対しても最適であり, 逆も成り立つ.

証明 $V(\Delta_1)$ は S に関して convex であること (cf. Schäl [1979]) を用いて, Benson [1979] の Theorem 4.2 を適用する \square

定常政策 f^∞ に対して

$$R_f(\Delta) = E_{f^\infty} \left[\sum_{i=1}^{\infty} \beta^{i-1} r(\Delta_i, a_i) \mid \Delta \right]$$

$$L_f(\Delta, a) = r(\Delta, a) + \beta \int_S R_f(\Delta') d\mathcal{G}(\Delta' \mid \Delta, a)$$

$$H_c = \{x \in \mathbb{R}^p \mid \langle c, x \rangle \leq 0\}, \quad c \in \mathbb{R}^p$$

$P_\pi(\Delta_1)$: 政策 π によって生ずる $A \times (S \times A)^\infty$ 上の確率測度

と定義する.

定理 3.1 定常政策 $(f^*)^\infty$ が VMDP に対して最適ならば, 各 $\Delta_1 \in S$ に対して,

$$L_{f^*}(\Delta_n, a_n) - R_{f^*}(\Delta_n) \in H_{c(\Delta_1)}, \quad P_{(f^*)^\infty}(\Delta_1) - \text{a.s. } \Delta_n, a_n \in F(\Delta_n), n \geq 1 \quad (3.1)$$

が成り立つような $c(\Delta_1) \in (\text{int } K^*)$, $\Delta_1 \in S$ が存在する.

証明 命題 3.1 より, $(f^*)^\infty$ が最適ならば $(f^*)^\infty$ はある $c(\Delta_1) \in (\text{int } K^*)$, $\Delta_1 \in S$ をもつ $\wedge \text{DP}(c)$ に対して最適である. 次の条件付期待値を考える.

$$J(\pi)(h_n) = E_\pi \left[\sum_{i=n}^{\infty} \beta^{i-n} r(\Delta_i, \Delta_i, a_i) \mid h_n \right]$$

ここで, $E_\pi[\cdot \mid h_n]$ は $\bar{\pi}_n \bar{\mathcal{G}} \bar{\pi}_{n+1} \bar{\mathcal{G}} \cdots$ による条件付期待値を

表わす (cf. Hinderer [1970], p. 80 および Appendix 3) .

特に, $\pi = (f^*)^\infty$ ならば

$$J(\pi)(h_n) = J((f^*)^\infty)(\Delta_1, \Delta_n) = \langle C(\Delta_1), R_{f^*}(\Delta_n) \rangle \quad (3.2)$$

また, $\pi = \{f^*, \dots, f^*, f_n, f^*, \dots\}$ ならば

$$J((f^*)^\infty)(\Delta_1, \Delta_n) \geq J(\pi)(h_n), \quad P_{(f^*)^\infty}(\Delta_1) - \text{a.s. } h_n \quad (3.3)$$

が成り立ち, $a_n = f(\Delta_n)$ において右辺を書き改めると,

$$J((f^*)^\infty)(\Delta_1, \Delta_n) \geq r(\Delta_1, \Delta_n, a_n) + \beta \int_S J((f^*)^\infty)(\Delta_1, \Delta_n) d q(\Delta_{n+1} | \Delta_n, a_n),$$

$$P_{(f^*)^\infty}(\Delta_1) - \text{a.s. } \Delta_n, a_n \in F(\Delta_n) \quad (3.4)$$

(3.2) を使って (3.4) を書き改めると結果を得る \square

定理 3.2 各 $\Delta_1 \in S$ に対して

$$L_{f^*}(\Delta_n, a_n) - R_{f^*}(\Delta_n) \in H_{C(\Delta_1)}, \quad P_\pi(\Delta_1) - \text{a.s. } \Delta_n, a_n \in F(\Delta_n), \quad n \geq 1, \quad \pi \in \Pi \quad (3.5)$$

が成り立つような $C(\Delta_1) \in (\text{int } K^*)$, $\Delta_1 \in S$ が存在すれば、定常政策 $(f^*)^\infty$ は VMDP に対して最適である.

証明

$$r(\Delta_1, \Delta_n, a_n) = \langle C(\Delta_1), r(\Delta_n, a_n) \rangle$$

を利得関数にもつ VDP(C) を考える. $u(\Delta_1, \Delta_n) = J((f^*)^\infty)(\Delta_1, \Delta_n)$

とおくと, (3.5) より

$$u(\Delta_1, \Delta_n) \geq r(\Delta_1, \Delta_n, a_n) + \beta \int_S u(\Delta_1, \Delta_{n+1}) d\mathcal{G}(\Delta_{n+1} | \Delta_n, a_n),$$

$$P_\pi(\Delta_1) - \text{a.s. } \Delta_n, a_n \in F(\Delta_n), n \geq 1, \pi \in \Pi$$
(3.6)

一般に、条件付期待値の性質より

$$E_\pi \left[\sum_{i=1}^{\infty} \{ \beta^i u(\Delta_1, \Delta_{i+1}) - E_\pi [\beta^i u(\Delta_1, \Delta_{i+1}) | \bar{h}_i] \} \middle| \Delta_1 \right] = 0$$
(3.7)

が成り立つ。ただし、 $\bar{h}_i = (\Delta_1, a_1, \dots, \Delta_i, a_i)$ 。ここで、(3.6)より

$$\begin{aligned} & E_\pi [\beta^i u(\Delta_1, \Delta_{i+1}) | \bar{h}_i] \\ &= \beta^i \int_S u(\Delta_1, \Delta_{i+1}) d\mathcal{G}(\Delta_{i+1} | \Delta_i, a_i) \\ &= \beta^{i-1} \left\{ r(\Delta_1, \Delta_i, a_i) + \beta \int_S u(\Delta_1, \Delta_{i+1}) d\mathcal{G}(\Delta_{i+1} | \Delta_i, a_i) \right\} \\ &\quad - \beta^{i-1} r(\Delta_1, \Delta_i, a_i) \\ &\leq \beta^{i-1} u(\Delta_1, \Delta_i) - \beta^{i-1} r(\Delta_1, \Delta_i, a_i). \end{aligned}$$

これを (3.7) に代入して、 $N \rightarrow \infty$ とすると、

$$J((\pi^*)^\infty)(\Delta_1) = u(\Delta_1, \Delta_1) \geq E_\pi \left[\sum_{i=1}^{\infty} \beta^{i-1} r(\Delta_1, \Delta_i, a_i) \middle| \Delta_1 \right] = J(\pi)(\Delta_1),$$

$\Delta_1 \in S.$

π は任意なので、 $(\pi^*)^\infty$ は $\text{NDP}(C)$ に対して最適である。したがって、 VMDP に対しても最適である \square

§ 4. 例

次のような VMDP を考える

$$K = R_+^2 = \{(x, y) \mid x \geq 0, y \geq 0\}$$

$$S = \{1, 2, 3\}, A = \{1, 2, 3\}, F(1) = F(2) = \{1, 2\}, F(3) = \{1\}$$

$$g(\{1\} \mid 1, 1) = g(\{1\} \mid 1, 2) = 1, g(\{2\} \mid 2, 1) = g(\{3\} \mid 2, 2) = 1,$$

$$g(\{3\} \mid 3, 1) = 1$$

$$r(1, 1) = (2, 2), r(1, 2) = (1, 4), r(2, 1) = (2, 2), r(2, 2) = (3, 1),$$

$$r(3, 1) = (3, 1)$$

この VMDP において, 定常政策 $\gamma^\infty: \gamma(1) = \gamma(2) = \gamma(3) = 1$ は最適であることを定理 3.2 を用いて示す.

$$I(\gamma^\infty)(1) = I(\gamma^\infty)(2) = (2/1-\beta, 2/1-\beta)$$

$$I(\gamma^\infty)(3) = (3/1-\beta, 1/1-\beta)$$

各 $a_i \in S$ に対して

$$D_{a_i} = \{ L_\gamma(a_n, a_n) - R_\gamma(a_n) \mid P_\pi(a_i)(a_n) > 0, a_n \in F(a_i), \pi \in \Pi \}$$

を求める.

$$D_1 : \begin{cases} r(1, 1) + \beta I(\gamma^\infty)(1) - I(\gamma^\infty)(1) = (0, 0) \\ r(1, 2) + \beta I(\gamma^\infty)(1) - I(\gamma^\infty)(1) = (-1, 2) \end{cases}$$

$$D_2 : \begin{cases} r(2, 1) + \beta I(\gamma^\infty)(2) - I(\gamma^\infty)(2) = (0, 0) \\ r(2, 2) + \beta I(\gamma^\infty)(2) - I(\gamma^\infty)(2) = (1, -1) \\ r(3, 1) + \beta I(\gamma^\infty)(3) - I(\gamma^\infty)(3) = (0, 0) \end{cases}$$

$$D_3 : r(3,1) + \beta I(\gamma^\infty)(3) - I(\gamma^\infty)(3) = (0,0)$$

図1より, 各 $\Delta_i \in S$ に対して (3.5) を満たすベクトル $C_{\Delta_i} \in (\text{int } K^*)$ が存在することがわかる. 故に, γ^∞ は最適である.

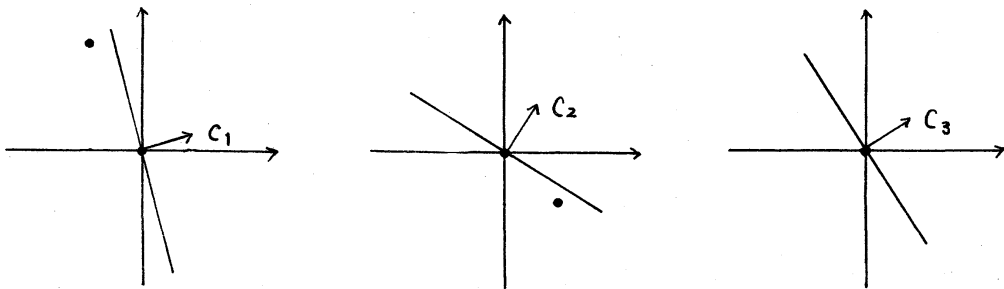


図 1

注意 スカラー化という方法で VMDP を考えるとき, 通常はベクトル $C(\Delta_i)$ は初期状態には依存せず一定である. しかし, そのような方法では上述の例の定常政策 γ^∞ の最適性は判定されない.

参考文献

- [1] H. P. Benson, An improved definition of proper efficiency for vector maximization with respect to cones, J. Math. Anal. Appl. 71 (1979), 232-241
- [2] N. Furukawa, Characterization of optimal policies in vector-valued Markovian decision processes, Math. Oper. Res. 5 (1980), 271-279

- [3] K. Hinderer, *Foundations of Non-stationary Dynamic Programming with Discrete Time Parameter* (Springer-Verlag, Berlin, 1970)
- [4] M. Schäl, On dynamic programming and statistical decision theory, *Ann. Probab.* 7 (1979), 832-885
- [5] D. J. White, Multi-objective infinite-horizon discounted Markov decision processes, *J. Math. Anal. Appl.* 89 (1982), 639-687